
Regularized Laplacian Estimation and Fast Eigenvector Approximation

Patrick O. Perry

Information, Operations, and Management Sciences
 NYU Stern School of Business
 New York, NY 10012
 pperry@stern.nyu.edu

Michael W. Mahoney

Department of Mathematics
 Stanford University
 Stanford, CA 94305
 mmahoney@cs.stanford.edu

Abstract

Recently, Mahoney and Orecchia demonstrated that popular diffusion-based procedures to compute a quick *approximation* to the first nontrivial eigenvector of a data graph Laplacian *exactly* solve certain regularized Semi-Definite Programs (SDPs). In this paper, we extend that result by providing a statistical interpretation of their approximation procedure. Our interpretation will be analogous to the manner in which ℓ_2 -regularized or ℓ_1 -regularized ℓ_2 -regression (often called Ridge regression and Lasso regression, respectively) can be interpreted in terms of a Gaussian prior or a Laplace prior, respectively, on the coefficient vector of the regression problem. Our framework will imply that the solutions to the Mahoney-Orecchia regularized SDP can be interpreted as regularized estimates of the pseudoinverse of the graph Laplacian. Conversely, it will imply that the solution to this regularized estimation problem can be computed very quickly by running, *e.g.*, the fast diffusion-based PageRank procedure for computing an approximation to the first nontrivial eigenvector of the graph Laplacian. Empirical results are also provided to illustrate the manner in which approximate eigenvector computation *implicitly* performs statistical regularization, relative to running the corresponding exact algorithm.

1 Introduction

Approximation algorithms and heuristic approximations are commonly used to speed up the running time of algorithms in machine learning and data analysis. In some cases, the outputs of these approximate procedures are “better” than the output of the more expensive exact algorithms, in the sense that they lead to more robust results or more useful results for the downstream practitioner. Recently, Mahoney and Orecchia formalized these ideas in the context of computing the first nontrivial eigenvector of a graph Laplacian [1]. Recall that, given a graph G on n nodes or equivalently its $n \times n$ Laplacian matrix L , the top nontrivial eigenvector of the Laplacian *exactly* optimizes the Rayleigh quotient, subject to the usual constraints. This optimization problem can equivalently be expressed as a vector optimization program with the objective function $f(x) = x^T L x$, where x is an n -dimensional vector, or as a Semi-Definite Program (SDP) with objective function $F(X) = \text{Tr}(LX)$, where X is an $n \times n$ symmetric positive semi-definite matrix. This first nontrivial vector is, of course, of widespread interest in applications due to its usefulness for graph partitioning, image segmentation, data clustering, semi-supervised learning, etc. [2, 3, 4, 5, 6, 7].

In this context, Mahoney and Orecchia asked the question: do popular diffusion-based procedures—such as running the Heat Kernel or performing a Lazy Random Walk or computing the PageRank function—to compute a quick *approximation* to the first nontrivial eigenvector of L solve some other regularized version of the Rayleigh quotient objective function *exactly*? Understanding this algorithmic-statistical tradeoff is clearly of interest if one is interested in very large-scale applications, where performing statistical analysis to derive an objective and then calling a black box solver to optimize that objective exactly might be too expensive. Mahoney and Orecchia answered the above question in the affirmative, with the interesting twist that the regularization is on the SDP

formulation rather than the usual vector optimization problem. That is, these three diffusion-based procedures exactly optimize a regularized SDP with objective function $F(X) + \frac{1}{\eta}G(X)$, for some regularization function $G(\cdot)$ to be described below, subject to the usual constraints.

In this paper, we extend the Mahoney-Orecchia result by providing a statistical interpretation of their approximation procedure. Our interpretation will be analogous to the manner in which ℓ_2 -regularized or ℓ_1 -regularized ℓ_2 -regression (often called Ridge regression and Lasso regression, respectively) can be interpreted in terms of a Gaussian prior or a Laplace prior, respectively, on the coefficient vector of the regression problem. In more detail, we will set up a sampling model, whereby the graph Laplacian is interpreted as an observation from a random process; we will posit the existence of a “population Laplacian” driving the random process; and we will then define an estimation problem: find the inverse of the population Laplacian. We will show that the maximum a posteriori probability (MAP) estimate of the inverse of the population Laplacian leads to a regularized SDP, where the objective function $F(X) = \text{Tr}(LX)$ and where the role of the penalty function $G(\cdot)$ is to encode prior assumptions about the population Laplacian. In addition, we will show that when $G(\cdot)$ is the log-determinant function then the MAP estimate leads to the Mahoney-Orecchia regularized SDP corresponding to running the PageRank heuristic. Said another way, the solutions to the Mahoney-Orecchia regularized SDP can be interpreted as regularized estimates of the pseudoinverse of the graph Laplacian. Moreover, by Mahoney and Orecchia’s main result, the solution to this regularized SDP can be computed very quickly—rather than solving the SDP with a black-box solver and rather computing explicitly the pseudoinverse of the Laplacian, one can simply run the fast diffusion-based PageRank heuristic for computing an approximation to the first nontrivial eigenvector of the Laplacian L .

The next section describes some background. Section 3 then describes a statistical framework for graph estimation; and Section 4 describes prior assumptions that can be made on the population Laplacian. These two sections will shed light on the computational implications associated with these prior assumptions; but more importantly they will shed light on the implicit prior assumptions associated with making certain decisions to speed up computations. Then, Section 5 will provide an empirical evaluation, and Section 6 will provide a brief conclusion. Additional discussion is available in the Appendix.

2 Background on Laplacians and diffusion-based procedures

A weighted symmetric graph G is defined by a vertex set $V = \{1, \dots, n\}$, an edge set $E \subset V \times V$, and a weight function $w : E \rightarrow \mathbb{R}_+$, where w is assumed to be symmetric (i.e., $w(u, v) = w(v, u)$). In this case, one can construct a matrix, $L_0 \in \mathbb{R}^{V \times V}$, called the combinatorial Laplacian of G :

$$L_0(u, v) = \begin{cases} -w(u, v) & \text{when } u \neq v, \\ d(u) - w(u, u) & \text{otherwise,} \end{cases}$$

where $d(u) = \sum_v w(u, v)$ is called the degree of u . By construction, L_0 is positive semidefinite. Note that the all-ones vector, often denoted $\mathbf{1}$, is an eigenvector of L_0 with eigenvalue zero, i.e., $L_0 \mathbf{1} = 0$. For this reason, $\mathbf{1}$ is often called trivial eigenvector of L_0 . Letting D be a diagonal matrix with $D(u, u) = d(u)$, one can also define a normalized version of the Laplacian: $L = D^{-1/2} L_0 D^{-1/2}$. Unless explicitly stated otherwise, when we refer to the Laplacian of a graph, we will mean the normalized Laplacian.

In many situations, e.g., to perform spectral graph partitioning, one is interested in computing the first *nontrivial* eigenvector of a Laplacian. Typically, this vector is computed “exactly” by calling a black-box solver; but it could also be approximated with an iteration-based method (such as the Power Method or Lanczos Method) or by running a random walk-based or diffusion-based method to the asymptotic state. These random walk-based or diffusion-based methods assign positive and negative “charge” to the nodes, and then they let the distribution of charge evolve according to dynamics derived from the graph structure. Three canonical evolution dynamics are the following:

Heat Kernel. Here, the charge evolves according to the heat equation $\frac{\partial H_t}{\partial t} = -LH_t$. Thus, the vector of charges evolves as $H_t = \exp(-tL) = \sum_{k=0}^{\infty} \frac{(-t)^k}{k!} L^k$, where $t \geq 0$ is a time parameter, times an input seed distribution vector.

PageRank. Here, the charge at a node evolves by either moving to a neighbor of the current node or teleporting to a random node. More formally, the vector of charges evolves as

$$R_\gamma = \gamma (I - (1 - \gamma) M)^{-1}, \quad (1)$$

where M is the natural random walk transition matrix associated with the graph and where $\gamma \in (0, 1)$ is the so-called teleportation parameter, times an input seed vector.

Lazy Random Walk. Here, the charge either stays at the current node or moves to a neighbor. Thus, if M is the natural random walk transition matrix associated with the graph, then the vector of charges evolves as some power of $W_\alpha = \alpha I + (1 - \alpha)M$, where $\alpha \in (0, 1)$ represents the “holding probability,” times an input seed vector.

In each of these cases, there is a parameter (t , γ , and the number of steps of the Lazy Random Walk) that controls the “aggressiveness” of the dynamics and thus how quickly the diffusive process equilibrates; and there is an input “seed” distribution vector. Thus, *e.g.*, if one is interested in global spectral graph partitioning, then this seed vector could be a vector with entries drawn from $\{-1, +1\}$ uniformly at random, while if one is interested in local spectral graph partitioning [8, 9, 10, 11], then this vector could be the indicator vector of a small “seed set” of nodes. See Appendix A for a brief discussion of local and global spectral partitioning in this context.

Mahoney and Orecchia showed that these three dynamics arise as solutions to SDPs of the form

$$\begin{aligned} & \underset{X}{\text{minimize}} && \text{Tr}(LX) + \frac{1}{\eta}G(X) \\ & \text{subject to} && X \succeq 0, \\ & && \text{Tr}(X) = 1, \\ & && XD^{1/2}1 = 0, \end{aligned} \tag{2}$$

where G is a penalty function (shown to be the generalized entropy, the log-determinant, and a certain matrix- p -norm, respectively [1]) and where η is a parameter related to the aggressiveness of the diffusive process [1]. Conversely, solutions to the regularized SDP of (2) for appropriate values of η can be computed *exactly* by running one of the above three diffusion-based procedures. Notably, when $G = 0$, the solution to the SDP of (2) is uu' , where u is the smallest nontrivial eigenvector of L . More generally and in this precise sense, the Heat Kernel, PageRank, and Lazy Random Walk dynamics can be seen as “regularized” versions of spectral clustering and Laplacian eigenvector computation. Intuitively, the function $G(\cdot)$ is acting as a penalty function, in a manner analogous to the ℓ_2 or ℓ_1 penalty in Ridge regression or Lasso regression, and by running one of these three dynamics one is *implicitly* making assumptions about the form of $G(\cdot)$. In this paper, we provide a statistical framework to make that intuition precise.

3 A statistical framework for regularized graph estimation

Here, we will lay out a simple Bayesian framework for estimating a graph Laplacian. Importantly, this framework will allow for regularization by incorporating prior information.

3.1 Analogy with regularized linear regression

It will be helpful to keep in mind the Bayesian interpretation of regularized linear regression. In that context, we observe n predictor-response pairs in $\mathbb{R}^p \times \mathbb{R}$, denoted $(x_1, y_1), \dots, (x_n, y_n)$; the goal is to find a vector β such that $\beta'x_i \approx y_i$. Typically, we choose β by minimizing the residual sum of squares, *i.e.*, $F(\beta) = \text{RSS}(\beta) = \sum_i \|y_i - \beta'x_i\|_2^2$, or a penalized version of it. For Ridge regression, we minimize $F(\beta) + \lambda\|\beta\|_2^2$; while for Lasso regression, we minimize $F(\beta) + \lambda\|\beta\|_1$.

The additional terms in the optimization criteria (*i.e.*, $\lambda\|\beta\|_2^2$ and $\lambda\|\beta\|_1$) are called penalty functions; and adding a penalty function to the optimization criterion can often be interpreted as incorporating prior information about β . For example, we can model y_1, \dots, y_n as independent random observations with distributions dependent on β . Specifically, we can suppose y_i is a Gaussian random variable with mean $\beta'x_i$ and known variance σ^2 . This induces a conditional density for the vector $y = (y_1, \dots, y_n)$:

$$p(y \mid \beta) \propto \exp\left\{-\frac{1}{2\sigma^2}F(\beta)\right\}, \tag{3}$$

where the constant of proportionality depends only on y and σ . Next, we can assume that β itself is random, drawn from a distribution with density $p(\beta)$. This distribution is called a prior, since it encodes prior knowledge about β . Without loss of generality, the prior density can be assumed to take the form

$$p(\beta) \propto \exp\{-U(\beta)\}. \tag{4}$$

Since the two random variables are dependent, upon observing y , we have information about β . This information is encoded in the posterior density, $p(\beta | y)$, computed via Bayes' rule as

$$p(\beta | y) \propto p(y | \beta) p(\beta) \propto \exp\{-\frac{1}{2\sigma^2} F(\beta) - U(\beta)\}. \quad (5)$$

The MAP estimate of β is the value that maximizes $p(\beta | y)$; equivalently, it is the value of β that minimizes $-\log p(\beta | y)$. In this framework, we can recover the solution to Ridge regression or Lasso regression by setting $U(\beta) = \frac{\lambda}{2\sigma^2} \|\beta\|_2^2$ or $U(\beta) = \frac{\lambda}{2\sigma^2} \|\beta\|_1$, respectively. Thus, Ridge regression can be interpreted as imposing a Gaussian prior on β , and Lasso regression can be interpreted as imposing a double-exponential prior on β .

3.2 Bayesian inference for the population Laplacian

For our problem, suppose that we have a connected graph with n nodes; or, equivalently, that we have L , the normalized Laplacian of that graph. We will view this observed graph Laplacian, L , as a “sample” Laplacian, *i.e.*, as random object whose distribution depends on a true “population” Laplacian, \mathcal{L} . As with the linear regression example, this induces a conditional density for L , to be denoted $p(L | \mathcal{L})$. Next, we can assume prior information about the population Laplacian in the form of a prior density, $p(\mathcal{L})$; and, given the observed Laplacian, we can estimate the population Laplacian by maximizing its posterior density, $p(\mathcal{L} | L)$.

Thus, to apply the Bayesian formalism, we need to specify the conditional density of L given \mathcal{L} . In the context of linear regression, we assumed that the observations followed a Gaussian distribution. A graph Laplacian is not just a single observation—it is a positive semidefinite matrix with a very specific structure. Thus, we will take L to be a random object with expectation \mathcal{L} , where \mathcal{L} is another normalized graph Laplacian. Although, in general, \mathcal{L} can be distinct from L , we will require that the nodes in the population and sample graphs have the same degrees. That is, if $d = (d(1), \dots, d(n))$ denotes the “degree vector” of the graph, and $D = \text{diag}(d(1), \dots, d(n))$, then we can define

$$\mathcal{X} = \{X : X \succeq 0, XD^{1/2}1 = 0, \text{rank}(X) = n - 1\}, \quad (6)$$

in which case the population Laplacian and the sample Laplacian will both be members of \mathcal{X} . To model L , we will choose a distribution for positive semi-definite matrices analogous to the Gaussian distribution: a scaled Wishart matrix with expectation \mathcal{L} . Note that, although it captures the trait that L is positive semi-definite, this distribution does not accurately model every feature of L . For example, a scaled Wishart matrix does not necessarily have ones along its diagonal. However, the mode of the density is at \mathcal{L} , a Laplacian; and for large values of the scale parameter, most of the mass will be on matrices close to \mathcal{L} . Appendix B provides a more detailed heuristic justification for the use of the Wishart distribution.

To be more precise, let $m \geq n - 1$ be a scale parameter, and suppose that L is distributed over \mathcal{X} as a $\frac{1}{m}$ Wishart(\mathcal{L}, m) random variable. Then, $\mathbb{E}[L | \mathcal{L}] = \mathcal{L}$, and L has conditional density

$$p(L | \mathcal{L}) \propto \frac{\exp\{-\frac{m}{2} \text{Tr}(L\mathcal{L}^+)\}}{|\mathcal{L}|^{m/2}}, \quad (7)$$

where $|\cdot|$ denotes pseudodeterminant (product of nonzero eigenvalues). The constant of proportionality depends only on L , d , m , and n ; and we emphasize that the density is supported on \mathcal{X} . Eqn. (7) is analogous to Eqn. (3) in the linear regression context, with $1/m$, the inverse of the sample size parameter, playing the role of the variance parameter σ^2 . Next, suppose we have know that \mathcal{L} is a random object drawn from a prior density $p(\mathcal{L})$. Without loss of generality,

$$p(\mathcal{L}) \propto \exp\{-U(\mathcal{L})\}, \quad (8)$$

for some function U , supported on a subset $\bar{\mathcal{X}} \subseteq \mathcal{X}$. Eqn. (8) is analogous to Eqn. (4) from the linear regression example. Upon observing L , the posterior distribution for \mathcal{L} is

$$p(\mathcal{L} | L) \propto p(L | \mathcal{L}) p(\mathcal{L}) \propto \exp\{-\frac{m}{2} \text{Tr}(L\mathcal{L}^+) + \frac{m}{2} \log |\mathcal{L}^+| - U(\mathcal{L})\}, \quad (9)$$

with support determined by $\bar{\mathcal{X}}$. Eqn. (9) is analogous to Eqn. (5) from the linear regression example. If we denote by $\hat{\mathcal{L}}$ the MAP estimate of \mathcal{L} , then it follows that $\hat{\mathcal{L}}^+$ is the solution to the program

$$\begin{aligned} & \underset{X}{\text{minimize}} && \text{Tr}(LX) + \frac{2}{m} U(X^+) - \log |X| \\ & \text{subject to} && X \in \bar{\mathcal{X}} \subseteq \mathcal{X}. \end{aligned} \quad (10)$$

Note the similarity with Mahoney-Orecchia regularized SDP of (2). In particular, if $\bar{\mathcal{X}} = \{X : \text{Tr}(X) = 1\} \cap \mathcal{X}$, then the two programs are identical except for the factor of $\log |X|$ in the optimization criterion.

4 A prior related to the PageRank procedure

Here, we will present a prior distribution for the population Laplacian that will allow us to leverage the estimation framework of Section 3; and we will show that the MAP estimate of \mathcal{L} for this prior is related to the PageRank procedure via the Mahoney-Orecchia regularized SDP. Appendix C presents priors that lead to the Heat Kernel and Lazy Random Walk in an analogous way; in both of these cases, however, the priors are data-dependent in the strong sense that they explicitly depend on the number of data points.

4.1 Prior density

The prior we will present will be based on neutrality and invariance conditions; and it will be supported on \mathcal{X} , *i.e.*, on the subset of positive-semidefinite matrices that was the support set for the conditional density defined in Eqn. (7). In particular, recall that, in addition to being positive semi-definite, every matrix in the support set has rank $n - 1$ and satisfies $XD^{1/2}1 = 0$. Note that because the prior depends on the data (via the orthogonality constraint induced by D), this is not a prior in the fully Bayesian sense; instead, the prior can be considered as part of an empirical or pseudo-Bayes estimation procedure.

The prior we will specify depends only on the eigenvalues of the normalized Laplacian, or equivalently on the eigenvalues of the pseudoinverse of the Laplacian. Let $\mathcal{L}^+ = \tau O \Lambda O'$ be the spectral decomposition of the pseudoinverse of the normalized Laplacian \mathcal{L} , where $\tau \geq 0$ is a scale factor, $O \in \mathbb{R}^{n \times n-1}$ is an orthogonal matrix, and $\Lambda = \text{diag}(\lambda(1), \dots, \lambda(n-1))$, where $\sum_v \lambda(v) = 1$. Note that the values $\lambda(1), \dots, \lambda(n-1)$ are unordered and that the vector $\lambda = (\lambda(1), \dots, \lambda(n-1))$ lies in the unit simplex. If we require that the distribution for λ be exchangeable (invariant under permutations) and neutral ($\lambda(v)$ independent of the vector $(\lambda(u)/(1 - \lambda(v)) : u \neq v)$, for all v), then the only non-degenerate possibility is that λ is Dirichlet-distributed with parameter vector (α, \dots, α) [12]. The parameter α , to which we refer as the “shape” parameter, must satisfy $\alpha > 0$ for the density to be defined. In this case,

$$p(\mathcal{L}) \propto p(\tau) \prod_{v=1}^{n-1} \lambda(v)^{\alpha-1}, \quad (11)$$

where $p(\tau)$ is a prior for τ . Thus, the prior weight on \mathcal{L} only depends on τ and Λ . One implication is that the prior is “nearly” rotationally invariant, in the sense that $p(P' \mathcal{L} P) = p(\mathcal{L})$ for any rank- $(n-1)$ projection matrix P satisfying $PD^{1/2}1 = 0$.

4.2 Posterior estimation and connection to PageRank

To analyze the MAP estimate associated with the prior of Eqn. (11) and to explain its connection with the PageRank dynamics, the following proposition is crucial.

Proposition 4.1. *Suppose the conditional likelihood for L given \mathcal{L} is as defined in (7) and the prior density for \mathcal{L} is as defined in (11). Define $\hat{\mathcal{L}}$ to be the MAP estimate of \mathcal{L} . Then, $[\text{Tr}(\hat{\mathcal{L}}^+)]^{-1} \hat{\mathcal{L}}^+$ solves the Mahoney-Orecchia regularized SDP (2), with $G(X) = -\log |X|$ and η as given in Eqn. (12) below.*

Proof. For \mathcal{L} in the support set of the posterior, define $\tau = \text{Tr}(\mathcal{L}^+)$ and $\Theta = \tau^{-1} \mathcal{L}^+$, so that $\text{Tr}(\Theta) = 1$. Further, $\text{rank}(\Theta) = n - 1$. Express the prior in the form of Eqn. (8) with function U given by

$$U(\mathcal{L}) = -\log\{p(\tau) |\Theta|^{\alpha-1}\} = -(\alpha-1) \log |\Theta| - \log p(\tau),$$

where, as before, $|\cdot|$ denotes pseudodeterminant. Using (9) and the relation $|\mathcal{L}^+| = \tau^{n-1} |\Theta|$, the posterior density for \mathcal{L} given L is

$$p(\mathcal{L} | L) \propto \exp \left\{ -\frac{m\tau}{2} \text{Tr}(L\Theta) + \frac{m+2(\alpha-1)}{2} \log |\Theta| + g(\tau) \right\},$$

where $g(\tau) = \frac{m(n-1)}{2} \log \tau + \log p(\tau)$. Suppose $\hat{\mathcal{L}}$ maximizes the posterior likelihood. Define $\hat{\tau} = \text{Tr}(\hat{\mathcal{L}}^+)$ and $\hat{\Theta} = [\hat{\tau}]^{-1} \hat{\mathcal{L}}^+$. In this case, $\hat{\Theta}$ must minimize the quantity $\text{Tr}(L\hat{\Theta}) - \frac{1}{\eta} \log |\hat{\Theta}|$, where

$$\eta = \frac{m\hat{\tau}}{m + 2(\alpha - 1)}. \quad (12)$$

Thus $\hat{\Theta}$ solves the regularized SDP (2) with $G(X) = -\log |X|$. \square

Mahoney and Orecchia showed that the solution to (2) with $G(X) = -\log |X|$ is closely related to the PageRank matrix, R_γ , defined in Eqn. (1). By combining Proposition 4.1 with their result, we get that the MAP estimate of \mathcal{L} satisfies $\hat{\mathcal{L}}^+ \propto D^{-1/2} R_\gamma D^{1/2}$; conversely, $R_\gamma \propto D^{1/2} \hat{\mathcal{L}}^+ D^{-1/2}$. Thus, the PageRank operator of Eqn. (1) can be viewed as a degree-scaled regularized estimate of the pseudoinverse of the Laplacian. Moreover, prior assumptions about the spectrum of the graph Laplacian have direct implications on the optimal teleportation parameter. Specifically Mahoney and Orecchia’s Lemma 2 shows how η is related to the teleportation parameter γ , and Eqn. (12) shows how the optimal η is related to prior assumptions about the Laplacian.

5 Empirical evaluation

In this section, we provide an empirical evaluation of the performance of the regularized Laplacian estimator, compared with the unregularized estimator. To do this, we need a ground truth population Laplacian \mathcal{L} and a noisily-observed sample Laplacian L . Thus, in Section 5.1, we construct a family of distributions for \mathcal{L} ; importantly, this family will be able to represent both low-dimensional graphs and expander-like graphs. Interestingly, the prior of Eqn. (11) captures some of the qualitative features of both of these types of graphs (as the shape parameter is varied). Then, in Section 5.2, we describe a sampling procedure for L which, superficially, has no relation to the scaled Wishart conditional density of Eqn. (7). Despite this model misspecification, the regularized estimator \hat{L}_η outperforms L for many choices of the regularization parameter η .

5.1 Ground truth generation and prior evaluation

The ground truth graphs we generate are motivated by the Watts-Strogatz “small-world” model [13]. To generate a ground truth population Laplacian, \mathcal{L} —equivalently, a population graph—we start with a two-dimensional lattice of width w and height h , and thus $n = wh$ nodes. Points in the lattice are connected to their four nearest neighbors, making adjustments as necessary at the boundary. We then perform s edge-swaps: for each swap, we choose two edges uniformly at random and then we swap the endpoints. For example, if we sample edges $i_1 \sim j_1$ and $i_2 \sim j_2$, then we replace these edges with $i_1 \sim j_2$ and $i_2 \sim j_1$. Thus, when $s = 0$, the graph is the original discretization of a low-dimensional space; and as s increases to infinity, the graph becomes more and more like a uniformly chosen 4-regular graph (which is an expander [14] and which bears similarities with an Erdős-Rényi random graph [15]). Indeed, each edge swap is a step of the Metropolis algorithm toward a uniformly chosen random graph with a fixed degree sequence. For the empirical evaluation presented here, $h = 7$ and $w = 6$; but the results are qualitatively similar for other values.

Figure 1 compares the expected order statistics (sorted values) for the Dirichlet prior of Eqn. (11) with the expected eigenvalues of $\Theta = \mathcal{L}^+ / \text{Tr}(\mathcal{L}^+)$ for the small-world model. In particular, in Figure 1(a), we show the behavior of the order statistics of a Dirichlet distribution on the $(n - 1)$ -dimensional simplex with scalar shape parameter α , as a function of α . For each value of the shape α , we generated a random $(n - 1)$ -dimensional Dirichlet vector, λ , with parameter vector (α, \dots, α) ; we computed the $n - 1$ order statistics of λ by sorting its components; and we repeated this procedure for 500 replicates and averaged the values. Figure 1(b) shows a corresponding plot for the ordered eigenvalues of Θ . For each value of s (normalized, here, by the number of edges μ , where $\mu = 2wh - w - h = 71$), we generated the normalized Laplacian, \mathcal{L} , corresponding to the random s -edge-swapped grid; we computed the $n - 1$ nonzero eigenvalues of Θ ; and we performed 1000 replicates of this procedure and averaged the resulting eigenvalues.

Interestingly, the behavior of the spectrum of the small-world model as the edge-swaps increase is qualitatively quite similar to the behavior of the Dirichlet prior order statistics as the shape parameter α increases. In particular, note that for small values of the shape parameter α the first few order-statistics are well-separated from the rest; and that as α increases, the order statistics become

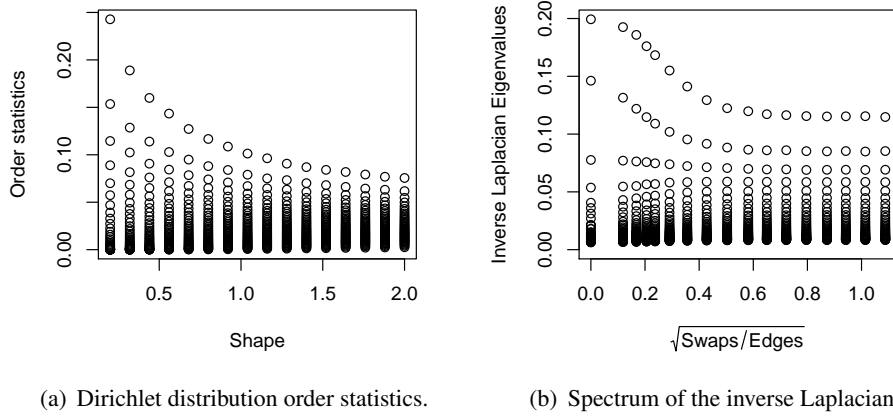


Figure 1: Analytical and empirical priors. 1(a) shows the Dirichlet distribution order statistics versus the shape parameter; and 1(b) shows the spectrum of Θ as a function of the rewiring parameter.

concentrated around $1/(n-1)$. Similarly, when the edge-swap parameter $s = 0$, the top two eigenvalues (corresponding to the width-wise and height-wise coordinates on the grid) are well-separated from the bulk; as s increases, the top eigenvalues quickly merge into the bulk; and eventually, as s goes to infinity, the distribution becomes very close that that of a uniformly chosen 4-regular graph.

5.2 Sampling procedure, estimation performance, and optimal regularization behavior

Finally, we evaluate the estimation performance of a regularized estimator of the graph Laplacian and compare it with an unregularized estimate. To do so, we construct the population graph \mathcal{G} and its Laplacian \mathcal{L} , for a given value of s , as described in Section 5.1. Let μ be the number of edges in \mathcal{G} . The sampling procedure used to generate the observed graph G and its Laplacian L is parameterized by the sample size m . (Note that this parameter is analogous to the Wishart scale parameter in Eqn. (7), but here we are sampling from a different distribution.) We randomly choose m edges with replacement from \mathcal{G} ; and we define sample graph G and corresponding Laplacian L by setting the weight of $i \sim j$ equal to the number of times we sampled that edge. Note that the sample graph G over-counts some edges in \mathcal{G} and misses others.

We then compute the regularized estimate $\hat{\mathcal{L}}_\eta$, up to a constant of proportionality, by solving (implicitly!) the Mahoney-Orecchia regularized SDP (2) with $G(X) = -\log |X|$. We define the unregularized estimate \hat{L} to be equal to the observed Laplacian, L . Given a population Laplacian \mathcal{L} , we define $\tau = \tau(\mathcal{L}) = \text{Tr}(\mathcal{L}^+)$ and $\Theta = \Theta(\mathcal{L}) = \tau^{-1}\mathcal{L}^+$. We define $\hat{\tau}_\eta$, $\hat{\tau}$, $\hat{\Theta}_\eta$, and $\hat{\Theta}$ similarly to the population quantities. Our performance criterion is the relative Frobenius error $\|\Theta - \hat{\Theta}_\eta\|_F / \|\Theta - \hat{\Theta}\|_F$, where $\|\cdot\|_F$ denotes the Frobenius norm ($\|A\|_F = [\text{Tr}(A'A)]^{1/2}$). Appendix D presents similar results when the performance criterion is the relative spectral norm error.

Figures 2(a), 2(b), and 2(c) show the regularization performance when $s = 4$ (an intermediate value) for three different values of m/μ . In each case, the mean error and one standard deviation around it are plotted as a function of $\eta/\bar{\tau}$, as computed from 100 replicates; here, $\bar{\tau}$ is the mean value of τ over all replicates. The implicit regularization clearly improves the performance of the estimator for a large range of η values. (Note that the regularization parameter in the regularized SDP (2) is $1/\eta$, and thus smaller values along the X-axis correspond to stronger regularization.) In particular, when the data are very noisy, *e.g.*, when $m/\mu = 0.2$, as in Figure 2(a), improved results are seen only for very strong regularization; for intermediate levels of noise, *e.g.*, $m/\mu = 1.0$, as in Figure 2(b), (in which case m is chosen such that G and \mathcal{G} have the same number of edges counting multiplicity), improved performance is seen for a wide range of values of η ; and for low levels of noise, Figure 2(c) illustrates that improved results are obtained for moderate levels of implicit regularization. Figures 2(d) and 2(e) illustrate similar results for $s = 0$ and $s = 32$.

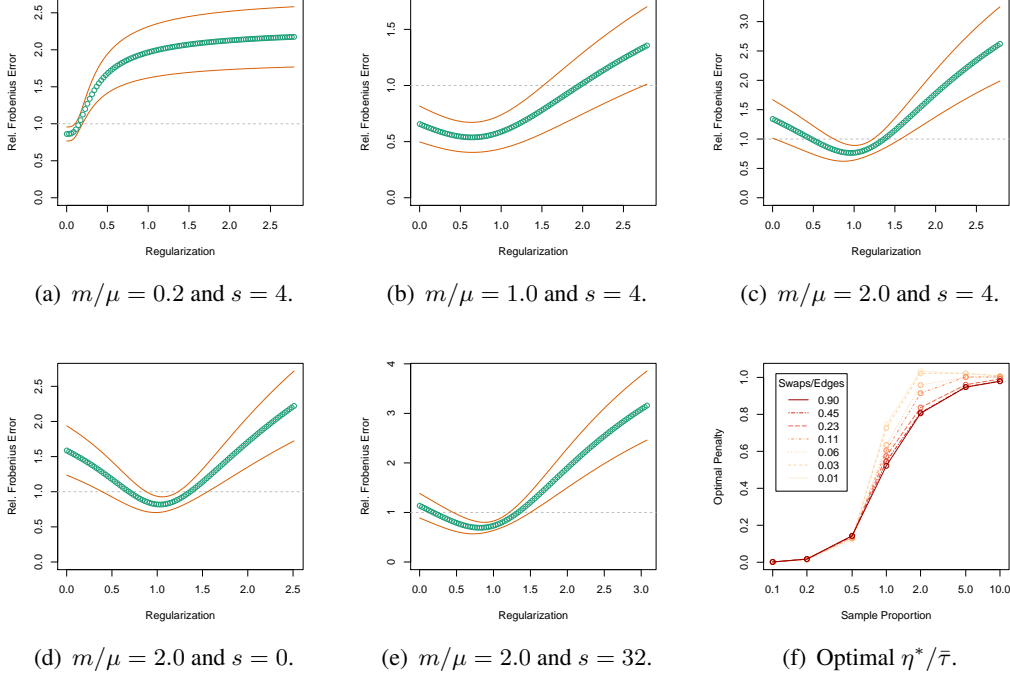


Figure 2: Regularization performance. 2(a) through 2(e) plot the relative Frobenius norm error, versus the (normalized) regularization parameter $\eta/\bar{\tau}$. Shown are plots for various values of the (normalized) number of edges, m/μ , and the edge-swap parameter, s . Recall that the regularization parameter in the regularized SDP (2) is $1/\eta$, and thus smaller values along the X-axis correspond to stronger regularization. 2(f) plots the optimal regularization parameter $\eta^*/\bar{\tau}$ as a function of sample proportion for different fractions of edge swaps.

As when regularization is implemented explicitly, in all these cases, we observe a “sweet spot” where there is an optimal value for the implicit regularization parameter. Figure 2(f) illustrates how the optimal choice of η depends on parameters defining the population Laplacians and sample Laplacians. In particular, it illustrates how η^* , the optimal value of η (normalized by $\bar{\tau}$), depends on the sampling proportion m/μ and the swaps per edges s/μ . Observe that as the sample size m increases, η^* converges monotonically to $\bar{\tau}$; and, further, that higher values of s (corresponding to more expander-like graphs) correspond to higher values of η^* . Both of these observations are in direct agreement with Eqn. (12).

6 Conclusion

We have provided a statistical interpretation for the observation that popular diffusion-based procedures to compute a quick approximation to the first nontrivial eigenvector of a data graph Laplacian exactly solve a certain regularized version of the problem. One might be tempted to view our results as “unfortunate,” in that it is not straightforward to interpret the priors presented in this paper. Instead, our results should be viewed as making explicit the implicit prior assumptions associated with making certain decisions (that are *already* made in practice) to speed up computations.

Several extensions suggest themselves. The most obvious might be to try to obtain Proposition 4.1 with a more natural or empirically-plausible model than the Wishart distribution; to extend the empirical evaluation to much larger and more realistic data sets; to apply our methodology to other widely-used approximation procedures; and to characterize when implicitly regularizing an eigenvector leads to better statistical behavior in downstream applications where that eigenvector is used. More generally, though, we expect that understanding the algorithmic-statistical tradeoffs that we have illustrated will become increasingly important in very large-scale data analysis applications.

References

- [1] M. W. Mahoney and L. Orecchia. Implementing regularization implicitly via approximate eigenvector computation. In *Proceedings of the 28th International Conference on Machine Learning*, pages 121–128, 2011.
- [2] D.A. Spielman and S.-H. Teng. Spectral partitioning works: Planar graphs and finite element meshes. In *FOCS '96: Proceedings of the 37th Annual IEEE Symposium on Foundations of Computer Science*, pages 96–107, 1996.
- [3] S. Guattery and G.L. Miller. On the quality of spectral separators. *SIAM Journal on Matrix Analysis and Applications*, 19:701–719, 1998.
- [4] J. Shi and J. Malik. Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):888–905, 2000.
- [5] M. Belkin and P. Niyogi. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural Computation*, 15(6):1373–1396, 2003.
- [6] T. Joachims. Transductive learning via spectral graph partitioning. In *Proceedings of the 20th International Conference on Machine Learning*, pages 290–297, 2003.
- [7] J. Leskovec, K.J. Lang, A. Dasgupta, and M.W. Mahoney. Community structure in large networks: Natural cluster sizes and the absence of large well-defined clusters. *Internet Mathematics*, 6(1):29–123, 2009. Also available at: arXiv:0810.1355.
- [8] D.A. Spielman and S.-H. Teng. Nearly-linear time algorithms for graph partitioning, graph sparsification, and solving linear systems. In *STOC '04: Proceedings of the 36th annual ACM Symposium on Theory of Computing*, pages 81–90, 2004.
- [9] R. Andersen, F.R.K. Chung, and K. Lang. Local graph partitioning using PageRank vectors. In *FOCS '06: Proceedings of the 47th Annual IEEE Symposium on Foundations of Computer Science*, pages 475–486, 2006.
- [10] F.R.K. Chung. The heat kernel as the pagerank of a graph. *Proceedings of the National Academy of Sciences of the United States of America*, 104(50):19735–19740, 2007.
- [11] M. W. Mahoney, L. Orecchia, and N. K. Vishnoi. A spectral algorithm for improving graph partitions with applications to exploring data graphs locally. Technical report. Preprint: arXiv:0912.0681 (2009).
- [12] J. Fabius. Two characterizations of the Dirichlet distribution. *The Annals of Statistics*, 1(3):583–587, 1973.
- [13] D.J. Watts and S.H. Strogatz. Collective dynamics of small-world networks. *Nature*, 393:440–442, 1998.
- [14] S. Hoory, N. Linial, and A. Wigderson. Expander graphs and their applications. *Bulletin of the American Mathematical Society*, 43:439–561, 2006.
- [15] B. Bollobas. *Random Graphs*. Academic Press, London, 1985.
- [16] W.E. Donath and A.J. Hoffman. Lower bounds for the partitioning of graphs. *IBM Journal of Research and Development*, 17:420–425, 1973.
- [17] M. Fiedler. A property of eigenvectors of nonnegative symmetric matrices and its application to graph theory. *Czechoslovak Mathematical Journal*, 25(100):619–633, 1975.
- [18] M. Mihail. Conductance and convergence of Markov chains—a combinatorial treatment of expanders. In *Proceedings of the 30th Annual IEEE Symposium on Foundations of Computer Science*, pages 526–531, 1989.
- [19] P. Yu. Chebotarev and E. V. Shamis. On proximity measures for graph vertices. *Automation and Remote Control*, 59(10):1443–1459, 1998.
- [20] A. K. Chandra, P. Raghavan, W. L. Ruzzo, R. Smolensky, and P. Tiwari. The electrical resistance of a graph captures its commute and cover times. In *Proceedings of the 21st Annual ACM Symposium on Theory of Computing*, pages 574–586, 1989.

A Relationship with local and global spectral graph partitioning

In this section, we briefly describe the connections between our results and global and local versions of spectral partitioning, which were a starting point for this work.

The idea of spectral clustering is to approximate the best partition of the vertices of a connected graph into two pieces by using the nontrivial eigenvectors of a Laplacian (either the combinatorial or the normalized Laplacian). This idea has had a long history [16, 17, 2, 3, 4]. The simplest version of spectral clustering involves computing the first nontrivial eigenvector (or another vector with Rayleigh quotient close to that of the first nontrivial eigenvector) of L , and “sweeping” over that vector. For example, one can take that vector, call it x , and, for a given threshold K , define the two sets of the partition as

$$C(x, K) = \{i : x_i \geq K\}, \quad \text{and} \\ \bar{C}(x, K) = \{i : x_i < K\}.$$

Typically, this vector x is computed by calling a black-box solver, but it could also be approximated with an iteration-based method (such as the Power Method or Lanczos Method) or a random walk-based method (such as running a diffusive procedure or PageRank-based procedure to the asymptotic state). Far from being a “heuristic,” this procedure provides a *global partition* that (via Cheeger’s inequality and for an appropriate choice of K) satisfies provable quality-of-approximation guarantees with respect to the combinatorial optimization problem of finding the best “conductance” partition in the entire graph [18, 2, 3].

Although spectral clustering reduces a graph to a single vector—the smallest nontrivial eigenvector of the graph’s Laplacian—and then clusters the nodes using the information in that vector, it is possible to obtain much more information about the graph by looking at more than one eigenvector of the Laplacian. In particular, the elements of the pseudoinverse of the combinatorial Laplacian, L_0^+ , give local (*i.e.*, node-specific) information about random walks on the graph. The reason is that the pseudoinverse L_0^+ of the Laplacian is closely related to random walks on the graph. See, *e.g.* [19] for details. For example, it is known that the quantity $L_0^+(u, u) + L_0^+(v, v) - L_0^+(u, v) - L_0^+(v, u)$ is proportional to the commute time, a symmetrized version of the length of time before a random walker started at node u reaches node v , whenever u and v are in the same connected component [20]. Similarly, the elements of the pseudoinverse of the *normalized* Laplacian give degree-scaled measures of proximity between the nodes of a graph. It is likely that $L^+(u, v)$ has a probabilistic interpretation in terms of random walks on the graph, along the lines of our methodology, but we are not aware of any such interpretation. From this perspective, given L^+ and a cutoff value, K , we can define a *local partition* around node u via $P_K(u) = \{v : L^+(u, v) > K\}$. (Note that if v is in $P_K(u)$, then u is in $P_K(v)$; in addition, if the graph is disconnected, then there exists a K such that u and v are in the same connected component iff $v \in P_K(u)$.) We call clustering procedures based on this idea *local spectral partitioning*.

Although the naïve way of performing this local spectral partitioning, *i.e.*, to compute L^+ explicitly, is prohibitive for anything but very small graphs, these ideas form the basis for very fast local spectral clustering methods that employ truncated diffusion-based procedures to compute localized vectors with which to partition. For example, this idea can be implemented by performing a diffusion-based procedure with an input seed distribution vector localized on a node u and then sweeping over the resulting vector. This idea was originally introduced in [8] as a diffusion-based operational procedure that was local in a very strong sense and that led to Cheeger-like bounds analogous to those obtained with the usual global spectral partitioning; and this was extended and improved by [9, 10]. In addition, an optimization perspective on this was provided by [11]. Although [11] is local in a weaker sense, it does obtain local Cheeger-like guarantees from an explicit locally-biased optimization problem, and it provides an optimization ansatz that may be interpreted as a “local eigenvector.” See [8, 9, 10, 11] for details. Understanding the relationship between the “operational procedure versus optimization ansatz” perspectives was the origin of [1] and thus of this work.

B Heuristic justification for the Wishart density

In this section, we describe a sampling procedure for L which, in a very crude sense, leads approximately to a conditional Wishart density for $p(L \mid \mathcal{L})$.

Let G be a graph with vertex set $V = \{1, 2, \dots, n\}$, edge set $E = V \times V$ equipped with the equivalence relation $(u, v) = (v, u)$. Let ω be an edge weight function, and let \mathcal{L}_0 and \mathcal{L} be

the corresponding combinatorial and normalized Laplacians. Let Δ be a diagonal matrix with $\Delta(u, u) = \sum_v \omega(u, v)$, so that $\mathcal{L} = \Delta^{-1/2} \mathcal{L}_0 \Delta^{-1/2}$. Suppose the weights are scaled such that $\sum_{(u,v) \in E} \omega(u, v) = 1$, and suppose further that $\Delta(u, u) > 0$. We refer to $\omega(u, v)$ as the population weight of edge (u, v) .

A simple model for the sample graph is as follows: we sample m edges from E , randomly chosen according to the population weight function. That is, we see edges $(u_1, v_1), (u_2, v_2), \dots, (u_m, v_m)$, where the edges are all drawn independently and identically such that the probability of seeing edge (u, v) is determined by ω :

$$\mathbb{P}_\omega\{(u_1, v_1) = (u, v)\} = \omega(u, v).$$

Note that we will likely see duplicate edges and not every edge with a positive weight will get sampled. Then, we construct a weight function from the sampled edges, called the sample weight function, w , defined such that

$$w(u, v) = \frac{1}{m} \sum_{i=1}^m 1\{(u_i, v_i) = (u, v)\},$$

where $1\{\cdot\}$ is an indicator vector. In turn, we construct a sample combinatorial Laplacian, L_0 , defined such that

$$L_0(u, v) = \begin{cases} \sum_w w(u, w) & \text{when } u = v, \\ -w(u, v) & \text{otherwise.} \end{cases}$$

Let D be a diagonal matrix such that $D(u, u) = \sum_v w(u, v)$, and define $L = D^{-1/2} L_0 D^{-1/2}$. Letting \mathbb{E}_ω denote expectation with respect to the probability law \mathbb{P}_ω , note that $\mathbb{E}_\omega[w(u, v)] = \omega(u, v)$, that $\mathbb{E}_\omega L_0 = \mathcal{L}_0$, and that $\mathbb{E}_\omega D = \Delta$. Moreover, the strong law of large numbers guarantees that as m increases, these three quantities converge almost surely to their expectations. Further, Slutsky's theorem guarantees that $\sqrt{m}(L - \mathcal{L})$ and $\sqrt{m}\Delta^{-1/2}(L_0 - \mathcal{L}_0)\Delta^{-1/2}$ converge in distribution to the same limit. We use this large-sample behavior to approximate the distribution of L by the distribution of $\Delta^{-1/2} L_0 \Delta^{-1/2}$. Put simply, we treat the degrees as known.

The distribution of L_0 is completely determined by the edge sampling scheme laid out above. However, the exact form for the density involves an intractable combinatorial sum. Thus, we appeal to a crude approximation for the conditional density. The approximation works as follows:

1. For $i = 1, \dots, m$, define $x_i \in \mathbb{R}^n$ such that

$$x_i(u) = \begin{cases} +s_i & \text{when } u = u_i, \\ -s_i & \text{when } u = v_i, \\ 0 & \text{otherwise,} \end{cases}$$

where $s_i \in \{-1, +1\}$ is chosen arbitrarily. Note that $L_0 = \frac{1}{m} \sum_{i=1}^m x_i x_i'$.

2. Take s_i to be random, equal to $+1$ or -1 with probability $\frac{1}{2}$. Approximate the distribution of x_i by the distribution of a multivariate normal random variable, \tilde{x}_i , such that x_i and \tilde{x}_i have the same first and second moments.
3. Approximate the distribution of L_0 by the distribution of \tilde{L}_0 , where $\tilde{L}_0 = \frac{1}{m} \sum_{i=1}^m \tilde{x}_i \tilde{x}_i'$.
4. Use the asymptotic expansion above to approximate the distribution of L by the distribution of $\Delta^{-1/2} \tilde{L}_0 \Delta^{-1/2}$.

The next two lemmas derive the distribution of \tilde{x}_i and \tilde{L}_0 in terms of \mathcal{L} , allowing us to get an approximation for $p(L \mid \mathcal{L})$.

Lemma B.1. *With x_i and \tilde{x}_i defined as above,*

$$\mathbb{E}_\omega[x_i] = \mathbb{E}_\omega[\tilde{x}_i] = 0,$$

and

$$\mathbb{E}_\omega[x_i x_i'] = \mathbb{E}_\omega[\tilde{x}_i \tilde{x}_i'] = \mathcal{L}_0.$$

Proof. The random variable \tilde{x}_i is defined to have the same first and second moments as x_i . The first moment vanishes since $s_i \stackrel{d}{=} -s_i$ implies that $x_i \stackrel{d}{=} -x_i$. For the second moments, note that when $u \neq v$,

$$\mathbb{E}_\omega[x_i(u) x_i(v)] = -s_i^2 \mathbb{P}_\omega\{(u_i, v_i) = (u, v)\} = -\omega(u, v) = \mathcal{L}_0(u, v).$$

Likewise,

$$\mathbb{E}_\omega[\{x_i(u)\}^2] = \sum_v \mathbb{P}_\omega\{(u_i, v_i) = (u, v)\} = \sum_v \omega(u, v) = \mathcal{L}_0(u, u). \quad \square$$

Lemma B.2. *The random matrix \tilde{L}_0 is distributed as $\frac{1}{m}$ Wishart(\mathcal{L}_0, m) random variable. This distribution is supported on the set of positive-semidefinite matrices with the same nullspace as \mathcal{L}_0 . When $m \geq \text{rank}(\mathcal{L}_0)$, the distribution has a density on this space given by*

$$f(\tilde{L}_0 \mid \mathcal{L}_0, m) \propto \frac{|\tilde{L}_0|^{(m-\text{rank}(\mathcal{L})-1)/2} \exp\{-\frac{m}{2} \text{Tr}(\tilde{L}_0 \mathcal{L}_0^+)\}}{|\mathcal{L}_0|^{m/2}} \quad (13)$$

where the constant of proportionality depends only on m and n and where $|\cdot|$ denotes pseudodeterminant (product of nonzero eigenvalues).

Proof. Since $m\tilde{L}$ is a sum of m outer products of multivariate Normal($0, \mathcal{L}_0$), it is Wishart distributed (by definition). Suppose $\text{rank}(\mathcal{L}_0) = r$ and $U \in \mathbb{R}^{n \times r}$ is a matrix whose columns are the eigenvectors of \mathcal{L}_0 . Note that $U' \tilde{x}_i \stackrel{d}{=} \text{Normal}(0, U' \mathcal{L}_0 U)$, and that $U' \mathcal{L}_0 U$ has full rank. Thus, $U' \tilde{L}_0 U$ has a density over the space of $r \times r$ positive-semidefinite matrices whenever $m \geq r$. The density of $U' \tilde{L}_0 U$ is exactly equal to $f(\tilde{L}_0 \mid \mathcal{L}_0, m)$, defined above. \square

Using the previous lemma, the random variable $\tilde{L} = \Delta^{-1/2} \tilde{L}_0 \Delta^{-1/2}$ has density

$$f(\tilde{L} \mid \mathcal{L}, m) \propto \frac{|\Delta^{1/2} \tilde{L} \Delta^{1/2}|^{(m-\text{rank}(\mathcal{L})-1)/2} \exp\{-\frac{m}{2} \text{Tr}(\Delta^{1/2} \tilde{L} \Delta^{1/2} \mathcal{L}_0^+)\}}{|\Delta^{1/2} \mathcal{L}_0 \Delta^{1/2}|^{m/2}},$$

where we have used that $\text{rank}(\mathcal{L}_0) = \text{rank}(\mathcal{L})$, and the constant of proportionality depends on m , n , $\text{rank}(\mathcal{L})$, and Δ . Then, if we approximate $|\Delta^{1/2} \tilde{L} \Delta^{1/2}| \approx |\Delta| |\tilde{L}|$ and $\Delta^{1/2} \mathcal{L}_0^+ \Delta^{1/2} \approx \mathcal{L}^+$, then f is “approximately” the density of a $\frac{1}{m}$ Wishart(\mathcal{L}, m) random variable. These last approximations are necessary because \tilde{L} and \mathcal{L}_0 are rank-degenerate.

To conclude, we do not want to overstate the validity of this heuristic justification. In particular, it makes three key approximations:

1. the true degree matrix Δ can be approximated by the observed degree matrix D ;
2. the distribution of x_i , a sparse vector, is well approximated \tilde{x}_i , a Gaussian (dense) vector;
3. the quantities $|\Delta^{1/2} \tilde{L} \Delta^{1/2}|$ and $\Delta^{1/2} \mathcal{L}_0^+ \Delta^{1/2}$ can be replaced with $|\Delta| |\tilde{L}|$ and \mathcal{L}^+ .

None of these approximations hold in general, though as argued above, the first is plausible if m is large relative to n . Likewise, since \tilde{L} and \mathcal{L} are nearly full rank, the third approximation is likely not too bad. The biggest leap of faith is the second approximation. Note, *e.g.*, that despite their first moments being equal, the second moments of $\tilde{x}_i \tilde{x}_i'$ and $x_i x_i'$ differ.

C Other priors and the relationship to Heat Kernel and Lazy Random Walk

There is a straightforward generalization of Proposition 4.1 to other priors. In this section, we state it, and we observe connections with the Heat Kernel and Lazy Random Walk procedures.

Proposition C.1. *Suppose the conditional likelihood for L given \mathcal{L} is as defined in (7) and the prior density for \mathcal{L} is of the form*

$$p(\mathcal{L}) \propto p(\tau) |\Theta|^{-m/2} \exp\{-q(\tau) G(\Theta)\}, \quad (14)$$

where $\tau = \text{Tr}(\mathcal{L}^+)$, $\Theta = \tau^{-1} \mathcal{L}^+$, and p and q are functions with $q(\tau) > 0$ over the support of the prior. Define $\hat{\mathcal{L}}$ to be the MAP estimate of \mathcal{L} . Then, $[\text{Tr}(\hat{\mathcal{L}}^+)]^{-1} \hat{\mathcal{L}}^+$ solves the Mahoney-Orecchia regularized SDP (2), with G the same as in the expression (14) for $p(\mathcal{L})$ and with

$$\eta = \frac{m \hat{\tau}}{2 q(\hat{\tau})},$$

where $\hat{\tau} = \text{Tr}(\hat{\mathcal{L}}^+)$.

The proof of this proposition is a straightforward generalization of the proof of Proposition 4.1 and is thus omitted. Note that we recover the result of Proposition 4.1 by setting $G(\Theta) = -\log |\Theta|$ and $q(\tau) = \frac{m}{2} + \alpha - 1$. In addition, by choosing $G(\cdot)$ to be the generalized entropy or the matrix p -norm penalty of [1], we obtain variants of the Mahoney-Orecchia regularized SDP (2) with the regularization term $G(\cdot)$. By then combining Proposition C.1 with their result, we get that the MAP estimate of \mathcal{L} is related to the Heat Kernel and Lazy Random Walk procedures, respectively, in a manner analogous to what we saw in Section 4 with the PageRank procedure. In both of these other cases, however, the prior $p(\mathcal{L})$ is data-dependent in the strong sense that it explicitly depends on the number of data points; and, in addition, the priors for these other cases do not correspond to any well-recognizable parametric distribution.

D Regularization performance with respect to the relative spectral error

In this section, we present Figure 3, which shows the regularization performance for our empirical evaluation, when the performance criterion is the relative spectral norm error, *i.e.*, $\|\Theta - \hat{\Theta}_\eta\|_2 / \|\Theta - \hat{\Theta}\|_2$, where $\|\cdot\|_2$ denotes spectral norm of a matrix (which is the largest singular value of that matrix). See Section 5.2 for details of the setup. Note that these results are very similar to those for the relative Frobenius norm error that are presented in Figure 2.

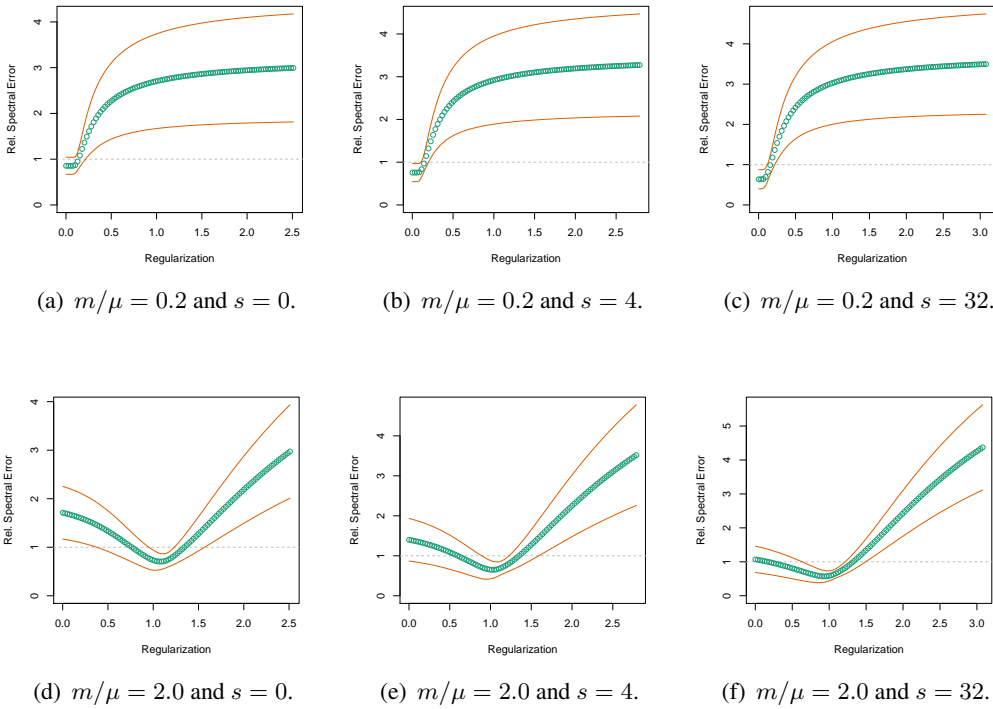


Figure 3: Regularization performance, as measured with the relative spectral norm error, versus the (normalized) regularization parameter $\eta/\bar{\tau}$. Shown are plots for various values of the (normalized) number of edges, m/μ , and the edge-swap parameter, s . Recall that the regularization parameter in the regularized SDP (2) is $1/\eta$, and thus smaller values along the X-axis correspond to stronger regularization.